

CLAIMS

What is claimed is:

- 1 1. A method of determining if a query document matches one or more
2 documents in a database, the method comprising:
3 generating a bit profile of the query document based on the number of
4 bits required to encode each of a plurality of rows of pixels in the
5 document; and
6 comparing the bit profile of the query document against bit profiles
7 associated with a first plurality of documents from the database to
8 determine if the query document matches one or more of the first
9 plurality of documents.
- 1 2. The method of claim 1 further comprising:
2 performing spectral analysis on the bit profile of the query document to
3 determine global statistics of the query document; and
4 comparing the global statistics of the query document against global
5 statistics associated with a second plurality of documents from the
6 database to identify the first plurality of documents.
- 1 3. The method of claim 2 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of a dominant line spacing in the query document.
- 1 4. The method of claim 2 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation

3 of a proportion of the query document that is text.

1 5. The method of claim 2 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of a location of text in the query document.

1 6. The method of claim 2 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of text concentration in the document, the estimation of text
4 concentration indicating a lengthwise measure of a proportion of the
5 query document that is text.

1 7. The method of claim 1 further comprising precomputing the bit files
2 associated with the first plurality of documents and storing the
3 precomputed bit profiles in the database.

1 8. The method of claim 1 wherein comparing the bit profile of the query
2 document against bit profiles associated with the first plurality of
3 documents comprises cross correlating the bit profile of the query
4 document against the bit profiles associated with the first plurality of
5 documents from the database.

1 9. The method of claim 8 wherein cross correlating the bit profile of the
2 query document against the bit profiles associated with the first plurality
3 of documents from the database comprises generating respective vector
4 products of the bit profile of the query document the bit profiles
5 associated with the first plurality of documents from the database.

1 10. The method of claim 9 wherein the query document is determined to
2 match one or more of the first plurality of documents for which the
3 respective vector product exceeds a threshold.

1 11. A method of determining if a query document matches one or more
2 documents in a database, the method comprising:
3 identifying up endpoints and down endpoints in the query document,
4 the up endpoints representing tops of features in the query
5 document and the down endpoints representing bottoms of
6 features in the query document;
7 generating a set of descriptors for the query document based on locations
8 of the up endpoints and the down endpoints; and
9 comparing the set of descriptors for the query document against
10 respective sets of descriptors associated with the one or more
11 documents in the database to determine if the query document
12 matches at least one of the one or more documents.

1 12. The method of claim 11 wherein generating a set of descriptors for the
2 query document based on locations of the up endpoints and the down
3 endpoints comprises
4 identifying text lines in the query document based on concentrations of
5 up endpoints and down endpoints along scanlines of the query
6 document; and
7 generating the set of descriptors based on distances between selected up
8 endpoints and selected down endpoints within the text lines in the

9 query document.

1 13. The method of claim 12 wherein identifying text lines in the document
2 based on concentrations of up endpoints and down endpoints along
3 scanlines of the document comprises:
4 determining the number of up endpoints and the number of down
5 endpoints that lie on each of the scanlines; and
6 identifying respective pairs of scanlines that have a local maximum
7 number of up endpoints and a local maximum number of down
8 endpoints as text lines.

1 14. The method of claim 11 wherein the query document is in a compressed
2 form in which respective runs of pixels are encoded in one of a plurality
3 of encoding modes, and wherein identification of the up endpoints and
4 down endpoints is unaffected by the encoding mode.

1 15. A method of determining if a query document matches one or more
2 documents in a database, the method comprising:
3 generating a bit profile of the query document based on the number of
4 bits required to encode each of a plurality of rows of pixels in the
5 query document;
6 comparing the bit profile of the query document against bit profiles
7 associated with a first plurality of documents from the database to
8 identify one or more candidate documents;
9 identifying endpoint features in the query document;
10 generating a set of descriptors for the query document based on locations

11 of the endpoint features; and
12 comparing the set of descriptors for the query document against
13 respective sets of descriptors for the one or more candidate
14 documents to determine if the query document matches at least
15 one of the one or more candidate documents.

1 16. The method of claim 15 further comprising
2 performing spectral analysis on the bit profile of the query document to
3 determine global statistics of the query document; and
4 comparing the global statistics of the query document against global
5 statistics associated with a second plurality of documents from the
6 database to identify the first plurality of documents, the first
7 plurality of documents being a subset of the second plurality of
8 documents.

1 17. The method of claim 16 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of at least one of a dominant line spacing in the query document, a
4 proportion of the query document that is text, a location of text in the
5 query document, and a text concentration.

1 18. The method of claim 15 wherein comparing the bit profile of the query
2 document against bit profiles associated with the first plurality of
3 documents comprises cross correlating the bit profile of the query
4 document against the bit profiles associated with the first plurality of
5 documents from the database.

1 19. A method of generating a set of descriptors for identifying a document,
2 the method comprising:
3 identifying up endpoints and down endpoints in the document, the up
4 endpoints representing tops of features in the document and the
5 down endpoints representing bottoms of features in the document;
6 identifying text lines in the document based on concentrations of up
7 endpoints and down endpoints along scanlines of the document;
8 and
9 generating a set of descriptors based on distances between selected up
10 endpoints and selected down endpoints in the concentrations of up
11 endpoints and down endpoints.

1 20. The method of claim 19 wherein identifying text lines in the document
2 based on concentrations of up endpoints and down endpoints along
3 scanlines of the document comprises:
4 determining the number of up endpoints and the number of down
5 endpoints that lie on each of the scanlines; and
6 identifying respective pairs of scanlines that have a local maximum
7 number of up endpoints and a local maximum number of down
8 endpoints as text lines.

1 21. The method of claim 19 wherein identifying text lines in the document
2 based on concentrations of up endpoints and down endpoints along
3 scanlines of the document comprises:
4 determining a dominant line spacing in the document;

5 determining the number of up endpoints and the number of down
6 endpoints that lie on each of the scanlines; and
7 identifying as text lines respective scanline pairs in which the
8 constituent scanlines are separated by a distance less than the
9 dominant line spacing and in which the constituent scanlines
10 respectively have a local maximum number of up endpoints and a
11 local maximum number of down endpoints as text lines.

1 22. The method of claim 21 wherein the dominant line spacing is
2 determined based on spectral analysis of locations of the endpoints in
3 the document.

1 23. The method of claim 19 further comprising generating a respective
2 endpoint profile for each of the scanlines, the endpoint profile
3 including a count of up endpoints identified on the scanline and a
4 count of down endpoints identified on the scanline, and wherein
5 identifying text lines based on concentrations of up endpoints and down
6 endpoints along scanlines of the document comprises reducing all but
7 local maximums of the counts of up endpoints and the counts of down
8 endpoints in respective endpoint profiles.

1 24. The method of claim 19 wherein identifying text lines based on
2 concentrations of up endpoints and down endpoints along scanlines of
3 the document comprises:
4 generating a count of up endpoints and a count of down endpoints for
5 each of the scanlines;
6 identifying a first scanline within a locality of scanlines that has the

7 highest count of up endpoints;
8 reducing the count of up endpoints associated with each scanline within
9 the locality of scanlines except the first scanline;
10 identifying a second scanline within the locality of scanlines that has the
11 highest count of down endpoints; and
12 reducing the count of down endpoints associated with each scanline
13 within the locality of scanlines except the second scanline.

1 25. The method of claim 24 wherein identifying the first scanline within
2 the locality of scanlines that has the highest count of up endpoints
3 comprises:
4 determining a dominant line spacing of the document; and
5 defining the locality of scanlines to be scanlines within a range greater
6 than the dominant line spacing but less than twice the dominant
7 line spacing.

1 26. The method of claim 19 wherein generating a set of descriptors based on
2 distances between selected up endpoints and selected down endpoints
3 comprises defining an ascender zone and a descender zone for each of
4 the text lines, the selected up endpoints being up endpoints in the
5 ascender zone and the selected down endpoints being down endpoints
6 in the descender zone.

1 27. The method of claim 26 wherein defining an ascender zone and a
2 descender zone for each of the text lines comprises:
3 defining a region above an x-height line of a first text line of the text

4 lines to be the ascender zone for the first text line; and
5 defining a region below the baseline of the first text line to be the
6 descender zone for the first text line.

1 28. The method of claim 27 wherein the ascender zone of the first text line
2 is bounded in part by the descender zone for the preceding text line.

1 29 The method of claim 19 wherein generating a set of descriptors based on
2 distances between selected up endpoints and selected down endpoints
3 comprises generating for a first text line of the text lines a first descriptor
4 that includes a plurality of distance measurements, each distance
5 measurement indicating a distance between a reference point and a
6 respective endpoint of the selected up endpoints and the selected down
7 endpoints.

1 30. The method of claim 29 wherein the reference point is one of the
2 selected up endpoints and the selected down endpoints.

1 31. The method of claim 29 wherein each distance measurement indicating
2 the distance between the reference point and the respective endpoint is
3 a relative distance to another endpoint of the selected up endpoints and
4 the selected down endpoints.

1 32. The method of claim 19 wherein the document is in a compressed form
2 in which respective runs of pixels are encoded in one of a plurality of
3 encoding modes, and wherein identification of the up endpoints and
4 down endpoints is unaffected by the encoding mode.

1 33. The method of claim 19 wherein the document has been compressed
2 using Group 4 compression.

1 34. A method of generating information that can be used to identify a
2 document, the method comprising:
3 generating a bit profile based on the number of bits required to encode
4 each of a plurality of rows of pixels in the document; and
5 performing spectral analysis on the bit profile to determine global
6 statistics of the document.

1 35. The method of claim 34 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of a dominant line spacing in the document.

1 36. The method of claim 35 wherein generating an estimation of a
2 dominant line spacing comprises generating a power spectrum density
3 from the bit profile and calculating the estimation of the dominant line
4 spacing from a peak value in the power spectrum density.

1 37. The method of claim 34 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of a proportion of the document that is text.

1 38. The method of claim 37 wherein generating an estimation of a
2 proportion of the document that is text comprises generating a power
3 spectrum density from the bit profile and calculating the estimation of

4 the proportion of the document based on an energy under a peak value
5 in the power spectrum density.

1 39. The method of claim 34 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of a location of text in the document.

1 40. The method of claim 39 wherein generating an estimation of a location
2 of text in the document comprises:
3 applying a bandpass filter to the bit profile to generate a text energy
4 profile; and
5 determining a centroid of the text energy profile to be the estimation of
6 the location of text in the document.

1 41. The method of claim 40 wherein applying a bandpass filter to the bit
2 profile comprises:
3 determining a dominant line spacing frequency of the document; and
4 selecting a center frequency of the bandpass filter based on the dominant
5 line spacing frequency.

1 42. The method of claim 34 wherein performing spectral analysis on the bit
2 profile to determine global statistics comprises generating an estimation
3 of text concentration in the document, the estimation of text
4 concentration indicating a lengthwise measure of a proportion of the
5 document that is text.

1 43. The method of claim 42 wherein generating an estimation of text
2 concentration in the document comprises:
3 applying a bandpass filter to the bit profile to generate a text energy
4 profile; and
5 determining the estimation of the text concentration based on a length
6 of the text energy profile.

1 44. An article of manufacture including one or more computer-readable
2 media that embody a program of instructions to configure a processing
3 system to determine if a query document matches one or more
4 documents in a database, wherein the program of instructions, when
5 executed by one or more processors in the processing system, causes the
6 one or more processors to:
7 generate a bit profile of the query document based on the number of bits
8 required to encode each of a plurality of rows of pixels in the
9 document; and
10 compare the bit profile of the query document against bit profiles
11 associated with a first plurality of documents from the database to
12 determine if the query document matches one or more of the first
13 plurality of documents.

1 45. The article of claim 44 wherein the one or more computer-readable
2 media include one or more non-volatile storage devices.

1 46. The article of claim 44 wherein the one or more computer-readable

media include a propagated data signal.

47. An article of manufacture including one or more computer-readable media that embody a program of instructions to configure a processing system to determine if a query document matches one or more documents in a database, wherein the program of instructions, when executed by one or more processors in the processing system, causes the one or more processors to:

- identify up endpoints and down endpoints in the query document, the up endpoints representing tops of features in the query document and the down endpoints representing bottoms of features in the query document;
- generate a set of descriptors for the query document based on locations of the up endpoints and the down endpoints; and
- compare the set of descriptors for the query document against respective sets of descriptors associated with the one or more documents in the database to determine if the query document matches at least one of the one or more documents.

48. An article of manufacture including one or more computer-readable media that embody a program of instructions to configure a processing system to determine if a query document matches one or more documents in a database, wherein the program of instructions, when executed by one or more processors in the processing system, causes the one or more processors to:

- generate a bit profile of the query document based on the number of bits

8 required to encode each of a plurality of rows of pixels in the query
9 document;

10 compare the bit profile of the query document against bit profiles
11 associated with a first plurality of documents from the database to
12 identify one or more candidate documents;

13 identify endpoint features in the query document;

14 generate a set of descriptors for the query document based on locations
15 of the endpoint features; and

16 compare the set of descriptors for the query document against respective
17 sets of descriptors for the one or more candidate documents to
18 determine if the query document matches at least one of the one or
19 more candidate documents.

1 49. A data processing system comprising:

2 a database of document images; and

3 a computer that includes a processing unit and a memory, the memory
4 having stored therein a program of instructions to configure the
5 computer to determine if a query document matches one or more
6 documents in the database, wherein the program of instructions,
7 when executed by the processing unit of the computer, causes the
8 computer to:

9 generate a bit profile of the query document based on the number
10 of bits required to encode each of a plurality of rows of pixels
11 in the document; and

12 compare the bit profile of the query document against bit profiles
13 associated with a first plurality of documents from the

14 database to determine if the query document matches one or
15 more of the first plurality of documents.

1 50. A data processing system comprising:

2 a database of document images; and

3 a computer that includes a processing unit and a memory, the memory
4 having stored therein a program of instructions to configure the
5 computer to determine if a query document matches one or more
6 documents in the database, wherein the program of instructions,
7 when executed by the processing unit of the computer, causes the
8 computer to:

9 identify up endpoints and down endpoints in the query document,
10 the up endpoints representing tops of features in the query
11 document and the down endpoints representing bottoms of
12 features in the query document;

13 generate a set of descriptors for the query document based on

14 locations of the up endpoints and the down endpoints; and

15 compare the set of descriptors for the query document against

16 respective sets of descriptors associated with the one or more
17 documents in the database to determine if the query

18 document matches at least one of the one or more documents.

1 51. A data processing system comprising:

2 a database of document images; and

3 a computer that includes a processing unit and a memory, the memory

4 having stored therein a program of instructions to configure the

5 computer to determine if a query document matches one or more
6 documents in the database, wherein the program of instructions,
7 when executed by the processing unit of the computer, causes the
8 computer to:

9 generate a bit profile of the query document based on the number
10 of bits required to encode each of a plurality of rows of pixels
11 in the query document;

12 compare the bit profile of the query document against bit profiles
13 associated with a first plurality of documents from the
14 database to identify one or more candidate documents;

15 identify endpoint features in the query document;

16 generate a set of descriptors for the query document based on
17 locations of the endpoint features; and

18 compare the set of descriptors for the query document against
19 respective sets of descriptors for the one or more candidate
20 documents to determine if the query document matches at
21 least one of the one or more candidate documents.